

Transactions of NAS of Azerbaijan, (2010), vol. XXX, No 4, pp. 141-152.

Received February 02, 2010; Revised May 11, 2010.

KILLED MARKOV DECISION PROCESSES ON FINITE TIME INTERVAL FOR COUNTABLE MODELS¹

Nestor R. PAROLYA and Yaroslav I. YELEYKO

Abstract

We consider killed Markov decision processes for countable models on a finite time-interval. Existence of a uniform ε -optimal policy is proven. We show the correctness of the fundamental equation. The optimal control problem is reduced to a similar problem for the derived model. We receive an optimality equation and a method for the construction of simple optimal policies. The sufficiency of simple policies for countable models is proven. We show the correctness of the Markovian property. Additionally, a dynamic programming principle is considered.

Classification: 90C40.

Keywords: Markov decision process; correctness; optimality equation; uniform ε -optimal policy.

1. Introduction. Markov decision processes arise in the different areas of the economics, in particular for the economic work planning of the separate business, economic sector or entire economics. At the beginning of each period we can build a plan for the next period knowing the last achieved state. The system development can be described mathematically as a deterministic process if we assume that the position of the system at the end of each period is uniquely defined by the state at the end of the period and by a plan for this period.

It is necessary to consider the influence of such factors as meteorological conditions, demographic transition, demand fluctuations, the imperfection of the compound production processes coordination, scientific discoveries and inventions etc. Stochastic models take into account these factors: if we know the state at the beginning of the period and the plan, we can only calculate the probability distribution for the next period. Therefore, leaving aside the system states in the past periods we come to the idea of Markov decision process ("the future depends not on the past, but only on the present").

The Markov decision processes are well described in [1]: the definition of Markov decision process is given, the concept of "model" Z^μ is presented, the definition of policy π is given, the assessment of policy - $\omega(\pi)$ and ν - assessment of process Z^μ are defined, the existence of a uniform ε -optimal policy is proved, the optimality equation and method for simple optimal policies constructing are presented, the sufficient of simple policies for countable models is proved, the correctness of the Markovian property is shown and dynamic programming principle is considered.

In [1] the model does not take into account one risk factor, namely the probability of bankruptcy at some determined moment of time. As a result, we come to the idea of killed Markov decision process where the business can crash with some nonzero probability at every moment of time, with the exception of the initial state.

The concept of the killed Markov decision process brings us closer to the real economic system which is not common without risk.

¹ Revised and corrected version of the paper published in *Transactions of NAS of Azerbaijan*, (2010), vol. XXX, No 4, pp. 141-152.

2. Killed Markov decision process. Let $X_t (t = m, \dots, n)$ and let $A_t (t = m+1, \dots, n)$ be countable or finite sets and at least one of them is countable. To the arbitrary $a \in A_t$ is assigned a probability distribution $p(\cdot|a) = \mathbb{P}(x_t = x|a_t = a, x_{t-1})$ on X_t .

Definition 1. The function p which defines the law of the transition from A_t to X_t is called the **transition function**.

Definition 2. The point $x^* = x_m \in X_t$ is called **killed state**, and $p(x^*|a)$ - the **probability of kill** if $\mathbb{P}(x_{t+1} = x^*|a_t = a) = \mathbb{P}(x_{t+1} = x_m|a_t = a) \equiv p(x^*|a), x_m \in X_m$.

Remark 1. In other words, the system moves into the initial(home) state when it hits a killed state(process is killed).

From the definition of the killed state it follows:

$$\forall a \in A_t \exists x^* \in X_t : p(x^*|a) = 1 - \sum_{x \in X_t \setminus x^*} p(x|a) > 0.$$

Definition 3 (Killed Markov decision process). A killed Markov decision process on a time interval $[m, n]$ is defined through the following objects:

1. Sets X_m, \dots, X_n (spaces of states);
2. Sets A_{m+1}, \dots, A_n (spaces of actions);
3. The projection mapping $j : A \rightarrow X$ where $A = \bigcup_{t=m+1}^n A_t$, $X = \bigcup_{t=m}^n X_t$: $j(A_t) = X_{t-1} \setminus \{x^*\}, x^* \in X_{t-1}, (t = m+2, \dots, n)$ and $j(A_{m+1}) = X_m$;
4. The probability distribution $p(\cdot|a) = \mathbb{P}(x_t = x|a_t = a, x_{t-1})$ on X_t with killed states

$$\mathbb{P}(x_{t+1} = x^*|a_t = a) = \mathbb{P}(x_{t+1} = x_m|a_t = a) \equiv p(x^*|a) > 0;$$

5. The function q on A (reward function);
6. The function r on X_n (terminal reward);
7. The function c (crash function), defined on the killed states $c(x^*) = - \sum_{i=m+1}^t \max_{a_i \in A_i} q(a_i), x^* \in X_t, t = m+1, \dots, n$ (function c ensures a total bankruptcy - total loss of accumulated capital or more);
8. The initial distribution μ on X_m .

A stochastic process defined through (1-8) is called the **killed Markov decision process** or the **model** and it is denoted by Z_μ^* . If the initial distribution μ is concentrated at the point x , we shall write Z_x^* .

Definition 4. The trajectory $l = x_m a_{m+1} x_{m+1} \dots a_n x_n$ is called the **way**. The set of all ways we denote $L = X \times (X \times A)^n$.

Our goal is to find a decision method which maximizes the mathematical expectation of the assessment of way l :

$$I(l, x^*) = \sum_{t=m+1}^n [q(a_t) + c(x_t^*)] + r(x_n), \quad (2.1)$$

where:

$x^* = (x_{m+1}^*, \dots, x_n^*)$ - vector of killed states;

$l = x_m a_{m+1}, \dots, a_n x_n$ - way.

The decision method is meant to be some *policy*.

3. Policies.

Definition 5. Let $A(x) \subset A$ is the set of all available actions at the state $x \in X$. $\varphi(x) : X \rightarrow A(x)$ is called the **simple policy** if $\varphi(x_{t-1}) = a_t$ for arbitrary x_t which is not a killed state with the probability distribution $p(\cdot|a_t)(m < t \leq n)$ and x_m with the initial distribution μ .

Remark 2. When we use the simple policy $\varphi(x)$ we get the way $l = x_m a_{m+1}, \dots, a_n x_n$.

Definition 6. The mapping $\pi : H \rightarrow \pi(\cdot|h \in H)$ is called a **killed policy**, where $\pi(\cdot|h \in H)$ is a probability distribution on $A(x_{t-1})$ and $H = X \times (A \times X)^{t-1}$ is a space of histories up to epoch $m \leq t-1 \leq n$ ($h \in H \Leftrightarrow h = x_m a_{m+1}, \dots, a_{t-1} x_{t-1}$).

Remark 3. Obviously, $x_{t-1} \neq x^*$.

Definition 7. Killed policy $\pi(\cdot|h)$ is called a **Markov policy** if $\pi(\cdot|h) = \pi(\cdot|x_{t-1})$.

The next conceptions can not be well-defined without the assumption:

Assumption 1. The reward function q and the terminal reward function r have the **supremum**, $\exists \sup_{a \in A} q(a)$ and $\exists \sup_{x \in X_n} r(x)$.

Definition 8. Let $p(\cdot|a)$ be the transition function and let $\pi(\cdot|h)$ be a policy. Every initial distribution μ is assigned to a probability distribution P^* in the space L which has such the notation:

$$\begin{aligned} P^*(l, x^*) &= P^*(x_m a_{m+1}, \dots, a_n x_n, x_{m+1}^*, \dots, x_n^*) = \\ &= \mu(x_m) \pi(a_{m+1}|x_m) p(x_{m+1}|a_{m+1}) p(x_{m+1}^*|a_{m+1}) \cdot \dots \cdot \pi(a_n|h_{n-1}) p(x_n|a_n) p(x_n^*|a_n) \end{aligned} \quad (3.1)$$

Remark 4. After the definition of the measure P^* the way l can be interpreted as a stochastic process. Additionally, this process is called the Markov process if the policy π is a Markov policy.

For all functions ξ from space L the mathematical expectation of ξ is given by

$$E^*(\xi) = \sum_{l \in L} \xi(l) P^*(l, x^*) \quad (3.2)$$

The assessment (2.1) of the way l is an example of such function. Next, we denote its expectation ω :

$$\omega = E^* I(l, x^*) = E^* \left[\sum_{t=m+1}^n [q(a_t) + c(x_t^*)] + r(x_n) \right] \quad (3.3)$$

Definition 9 (Assessment of policy). The value ω from (3.3) is called the **assessment of policy** π and is the function of the variable π ($\omega = \omega(\pi)$) for the killed Markov decision process Z_μ^* .

The goal of the research is the maximization of function $\omega(\pi)$.

Definition 10 (Assessment of process). $\nu \equiv \sup_{\pi} \omega(\pi)$ is called the **assessment of killed Markov decision process** Z_μ^* or **assessment of initial distribution** μ .

Remark 5. $\nu(x^*) = c(x^*)$.

Definition 11 (ε -optimal policy). A killed policy π is called **ε -optimal** for Z_μ^* if $\forall \varepsilon > 0 : \omega(\mu, \pi) \geq \nu(\mu) - \varepsilon$.

Definition 12 (Uniform ε -optimal policy). A killed policy is called **uniform ε -optimal** or **ε -optimal for process** Z_μ^* if π is ε -optimal for Z_μ^* for all μ - initial distribution.

4. Existence of uniform ε -optimal policy. Let π_x is ε -optimal policy for process Z_x^* . Its existence follows from the definition of the supremum.

We want to build a killed policy π which is ε -optimal for the model Z^* by using a sequence of the killed policies π_x .

It's natural to use the policy π_x when x is a starting point. Formally,

$$\bar{\pi}(\cdot|h) = \pi_{x(h)}(\cdot|h) \quad (4.1)$$

where $x(h)$ - the initial state of history h . It is clear that formula (4.1) defines some policy $\bar{\pi}$ and this policy will be ε -optimal. It means that $\forall \varepsilon \geq 0 : \omega(x, \bar{\pi}) = \omega(x, \pi_x) \geq \nu(x) - \varepsilon, \forall x \in X_m$.

Proposition 1 (Existence of the uniform ε -optimal killed policy). *Every killed policy $\bar{\pi}$ from (4.1) which is ε -optimal, i.e.*

$$\omega(x, \bar{\pi}) \geq \nu(x) - \varepsilon, (x \in X_m), \forall \varepsilon \geq 0$$

is uniform ε -optimal. It means that $\forall \mu, \forall \varepsilon \geq 0 : \sup_{\pi} \omega(\mu, \pi) \leq \omega(\mu, \bar{\pi}) + \varepsilon$.

Proof. From (3.1)-(3.3) it follows that $\forall \pi$:

$$\omega(\mu, \pi) = \sum_{l \in L} I(l, x^*) P^*(l, x^*) = \sum_{X_m} \mu(x) \omega(x, \pi). \quad (4.2)$$

Hence, it appears

$$\omega(\mu, \pi) = \sum_{X_m} \mu(x) \omega(x, \pi) \leq \sum_{X_m} \mu(x) \nu(x) \leq \sum_{X_m} \mu(x) [\omega(x, \bar{\pi}) + \varepsilon] = \omega(\mu, \bar{\pi}) + \varepsilon.$$

From the received inequalities it follows that

$$\sup_{\pi} \omega(\mu, \pi) \leq \sum_{X_m} \mu(x) \nu(x), \quad (4.3)$$

$$\omega(\mu, \bar{\pi}) \geq \sum_{X_m} \mu(x) \nu(x) - \varepsilon. \quad (4.4)$$

According to the arbitrariness of $\varepsilon > 0$ we get now from (4.3) and (4.4)

$$\sup_{\pi} \omega(\mu, \pi) = \sum_{X_m} \mu(x) \nu(x) \leq \omega(\mu, \bar{\pi}) + \varepsilon. \quad (4.5)$$

So the policy $\bar{\pi}$ is uniform ε -optimal. **Proposition 1 is proved.**

Corollary 1. For all initial distributions μ :

$$\nu(\mu) = \mu \nu. \quad (4.6)$$

Proof. It follows from $\nu(\mu) = \sum_{X_m} \mu(x) \nu(x) = \mu \nu$.

Remark 6. Formulas (4.2) and (4.6) allow us to reduce the analysis of the processes Z_{μ}^* for all μ to the analysis of the processes Z_x^* , $\forall x \in X_m$.

The policy π is built of the sequence $\pi_x, (x \in X_m)$ and has the following property (1):

For all initial distribution of the state $x \in X_m$ the probability distributions in space L which are assigned to the policies π and π_x from (3.1) are equal.

Definition 13. If $\bar{\pi}$ satisfies the property (1) then $\bar{\pi}$ is called the **combination of policies** π_x .

5. Derived model and fundamental equation. The decision process is a quite number of consecutive steps. The first step is the choice of probability distribution on A_{m+1} which depends on initial state. Since the choice is taken every initial distribution μ on X_m accords with probability distribution $\dot{\mu}$ on X_{m+1} . Now we consider $\dot{\mu}$ as initial distribution in moment of time $m + 1$.

As a result, we divide our maximization problem by two problems:

1. Choose the optimal policy for the next moments of time for every initial distribution on X_{m+1} ;
2. Choose the first step according to maximum reward and maximum value of the optimal policy assessment in the next time moments for initial distribution $\dot{\mu}$.

Definition 14 (Derived model). *The model which is build of the model Z^* by deletion X_m and A_{m+1} is called the **derived model** and it is denoted \dot{Z}^* .*

Proposition 2 (Fundamental equation).

$$\omega(x, \pi) = \sum_{A(x)} \pi(a|x) \left(q(a) + \dot{\omega}(p_a, \pi_a) \right), \quad (5.1)$$

where $p_a = p(\cdot|a)$, $\pi_a(\cdot|\dot{h}) = \pi(\cdot|y a \dot{h})$,

$a \in A_{m+1}$, $y = j(a)$, \dot{h} is a history in model \dot{Z}^* .

The equation (5.1) is called **fundamental** and expresses the assessment ω of the random policy π in model Z^* in terms of the assessment $\dot{\omega}$ of some policies in the model \dot{Z}^* .

Proof. According to (4.2) we get

$$\dot{\omega}(p_a, \pi_a) = \sum_{X_{m+1}} p(y|a) \dot{\omega}(y, \pi_a) \quad (5.2)$$

Let consider the spaces of ways L and \dot{L} in the models Z^* and \dot{Z}^* . Let P^* is the probability distribution on L according to the initial state x and the policy π and let P_a^* is the probability distribution on \dot{L} according to the initial distribution p_a and the policy π_a .

According to (2.1) and (3.1) $\forall \dot{l} \in \dot{L}$ we get

$$I(x a \dot{l}, x^*) = q(a) + I(\dot{l}, x_{-1}^*) \quad (5.3)$$

$$P^*(x a \dot{l}, x^*) = \pi(a|x) P_a^*(\dot{l}, x_{-1}^*) \quad (5.4)$$

$$a \in A(x), x_{-1}^* = (x_{m+2}^*, \dots, x_n^*), (x_{m+1}^*, x_{-1}^*) = x^*.$$

Under the notations in (3.2) and (3.3) we get

$$\omega(x, \pi) = \sum_L P^*(l, x^*) I(l, x^*) \quad (5.5)$$

$$\dot{\omega}(p_a, \pi_a) = \sum_{\dot{L}} P_a^*(\dot{l}, x_{-1}^*) I(\dot{l}, x_{-1}^*) \quad (5.6)$$

The measure $P^*(l, x^*)$ is nonzero only for ways which have the starting point x , i.e., for $x a \dot{l}$. That is why by the substitution in (5.5) of the expression of $I(l, x^*)$ from (5.3) and the expression of $P^*(l, x^*)$ from (5.4), and according to (5.6) we get the fundamental equation (5.1). **Proposition 2 is proved.**

Remark 7. *The fundamental equation is correct even without **Assumption 1**.*

6. Reducing the problem of the optimal decision to analogical problem for the derived model. From fundamental equation (5.1) it follows the following inequality

$$\omega(x, \pi) \leq \sup_{A(x)} [q(a) + \dot{\omega}(p_a, \pi_a)] \leq \sup_{A(x)} [q(a) + \dot{\nu}(p_a)] \quad (6.1)$$

$\forall x \in X_m$ and for every π ($\dot{\nu}$ which is the assessment of model \dot{Z}^*).

We denote $u(a) = q(a) + \dot{\nu}(p_a)$, ($a \in A_{m+1}$) and call this value - **assessment of the action a** .

According to (4.3) and $\nu(x^*) = c(x^*)$ we get $u = U\dot{\nu}$ where operator U transforms functions on the non-killed states on X to the functions on A and is given by

$$Uf(a) = q(a) + \sum_y p(y|a)f(y) + \sum_{y^*} p(y^*|a)c(y^*) \quad (6.2)$$

where y and y^* are the non-killed states and the killed states, respectively.

Let the operator V transforms the functions on A into the functions on non-killed and non-terminal states on X and satisfies

$$Vg(x) = \sup_{a \in A(x)} g(a) \quad (6.3)$$

Let us write the inequality (6.1) by using the operator V :

$$\omega(x, \pi) \leq Vu(x).$$

Then we consider \sup_{π} of the right and the left part of $\omega(x, \pi) \leq Vu(x)$ and we get

$$\nu \leq Vu. \quad (6.4)$$

Remark 8. Later we show the conditions which assure the equality in (6.4).

Definition 15 (Product of policies). Let $\dot{\pi}$ be a killed policy in the model \dot{Z}^* and to $x \in X_m$ is assigned some probability distribution $\gamma(\cdot|x)$ on A_{m+1} which is concentrated on $A(x)$. When we choose on the first step an action a and on all other steps we use the killed policy $\dot{\pi}$ then we get the killed policy π in the model Z^* . This policy is called the **product of policies** γ and $\dot{\pi}$ and is denoted by $\gamma\dot{\pi}$. It has the expression

$$\pi(\cdot|h) = \begin{cases} \gamma(\cdot|x) & \text{for } h = x \in X_m, \\ \dot{\pi}(\cdot|h) & \text{for } h = xa\dot{h}. \end{cases}$$

Proposition 3. Let $\pi = \gamma\dot{\pi}$ is a product of the killed policies γ and $\dot{\pi}$. If $\dot{\pi}$ is uniform ε' -optimal for model \dot{Z}^* then:

$$\nu = Vu. \quad (6.4)$$

Proof. The fundamental equation (5.1) for a product of policies has the following expression

$$\omega(x, \gamma\dot{\pi}) = \sum_{A(x)} \gamma(a|x) (q(a) + \dot{\omega}(p_a, \dot{\pi})) \quad (6.5)$$

Since $\dot{\pi}$ is ε' -optimal (it exists $\forall \varepsilon' \geq 0$ according to Proposition 1.) we get $\dot{\omega}(p_a, \dot{\pi}) \geq \dot{\nu}(p_a) - \varepsilon'$, and according to appearance of u equation (6.5) transforms to

$$\omega(x, \gamma\dot{\pi}) \geq \sum_{A(x)} \gamma(a|x) u(a) - \varepsilon'.$$

Lets consider the set

$$A_\chi(x) = \{a : a \in A(x), u(a) \geq Vu(x) - \chi\} \quad (x \in X_m).$$

$A_\chi(x)$ is nonempty for all $\chi > 0$. Let $\gamma(\cdot|x)$ be a probability distribution on $A(x)$ which is concentrated on $A_\chi(x)$.

Then

$$\sum_{A(x)} \gamma(a|x)u(a) \geq Vu(x) - \chi.$$

Since $\varepsilon' + \chi \leq \varepsilon$ we get

$$\omega(x, \pi) \geq Vu(x) - \varepsilon, \quad (x \in X_m). \quad (6.6)$$

According to (6.4) and (6.6) **Proposition 3 is proved.**

Corollary 1. *The assessment ν of the model Z^* is expressed in terms of the assessment $\dot{\nu}$ of the model \dot{Z}^* in the following way:*

$$\nu = Vu, \quad u = U\dot{\nu} \quad (6.7)$$

where operators U and V are defined in (6.2) and (6.3);

Corollary 2. *For all $\chi > 0$ exists such $\psi(x) : X_m \rightarrow A_{m+1}(x)$:*

$$u(\psi(x)) \geq \nu(x) - \chi \quad (6.8)$$

Here $\gamma(\cdot|x)$ can be the distribution concentrated at one point $\psi(x) \in A_\chi(x)$.

Corollary 3. *Let ε' and χ be the arbitrary nonnegative numbers. If $\hat{\pi}$ is uniform ε' -optimal for the model \dot{Z}^* and ψ is such as in Corollary 3 then the killed policy $\psi\hat{\pi}$ is uniform $(\varepsilon' + \chi)$ -optimal for the model Z^* .*

7. Optimality equation. Method for the construction of simple optimal policies.

Let assume that in our model Z^* $m = 0$. Let consider the models $Z_0^*, Z_1^*, \dots, Z_n^*$ where $Z^* = Z_0^*$ and Z_t^* is a derived model of Z_{t-1}^* . Let denote the assessments ν and u of the model Z_t^* as ν_t and $u_{t+1}(\nu_t$ on X_t , u_{t+1} on A_{t+1}). The reward function q and the transition function p we denote q_t and p_t .

According to the results of section 6 we get

$$\nu_{t-1} = Vu_t, \quad u_t = U\nu_t \quad (1 \leq t \leq n) \quad (7.1)$$

where

$$U_t f(a) = q_t(a) + \sum_{y \in X_t} p_t(y|a)f(y) + p_t(y^*|a)c(y^*), \quad (a \in A_t, y^* \in X_t),$$

$$V_t g(x) = \sup_{A(x)} g(a), \quad (x \in X_{t-1}),$$

and $\nu_n = r$.

Equations (7.1) are called the **optimality equations**. Let $T_t = V_t U_t$ then the optimality equations transform to

$$\nu_{t-1} = T_t \nu_t. \quad (7.1')$$

From (7.1), (7.1') and the condition $\nu_n = r$ we calculate $\nu_n, \nu_{n-1}, \dots, \nu_0$. Then we choose the action $\psi_t(x) : X_{t-1} \rightarrow A_t(x)$ for which holds

$$u_t(\psi_t) \geq \nu_{t-1} - \chi_t. \quad (7.2)$$

$\forall t = 1, 2, \dots, n$ and for all nonnegative $\chi_1, \chi_2, \dots, \chi_n$.

According to *Corollary 3* of *Proposition 3* the simple policy $\varphi = \psi_1 \psi_2 \dots \psi_n$ is uniform ε -optimal for the model $Z^* = Z_0^*$ and $\varepsilon = \sum_{i=1}^n \chi_i$. The equation (7.2) can be rewritten as

$$T_{\psi_t} \nu_t \geq \nu_{t-1} - \chi_t, \quad (7.2')$$

where the operator T_{ψ_t} transforms functions on X_t to functions on X_{t-1} in the following way

$$T_{\psi_t} f(x) = q_t[\psi_t(x)] + \sum_{X_t} p(y|\psi_t(x))f(y) + p_t(y^*|a)c(y^*). \quad (7.3)$$

Proposition 4. *Let π be an arbitrary killed policy in the derived model Z_k^* ($k = 1, 2, \dots, n$) and let $\psi_t : X_{t-1} \rightarrow A_t(x)$ ($t = 1, 2, \dots, k$) are arbitrary too then*

$$\omega_0(x, \psi_1 \psi_2 \dots \psi_k \pi) = T_{\psi_1} T_{\psi_2} \dots T_{\psi_k} \omega_k(x, \pi), \quad (7.4)$$

Proof. It follows from the fundamental equation (5.1), formulas (5.2), (7.3) and the mathematical induction.

Remark 9. *It follows from (7.4): the result will not change if our decision process is killed at the moment of time k and the terminal reward as the assessment of policy π is taken.*

Remark 10. *If we can choose ψ_t with $\chi_t = 0$ in (7.2') $\forall t = 1..n$ then the simple policy $\varphi = \psi_1 \dots \psi_n$ is called uniform optimal.*

8. The sufficiency of the simple policies for countable models. The question arises: do we lose something by using only simple policies? The previous result can not give us the answer. It only makes our losses indefinitely small.

Theorem 1 (Sufficiency of the simple policies). *Let μ is a fixed initial distribution and let π is a arbitrary killed policy then there exists φ -simple policy such that*

$$\omega(\mu, \pi) \leq \omega(\mu, \varphi). \quad (8.1)$$

Proof. It follows from *Proposition 5* and *Proposition 6*.

Proposition 5. *For all μ and for all killed policies π there exists the Markov policy θ such that*

$$\omega(\mu, \theta) = \omega(\mu, \pi) \quad (8.2)$$

*These two policies are called **equivalent**.*

Proposition 6. *For all Markov policies θ there exists a simple policy φ such that*

$$\omega(\mu, \varphi) \geq \omega(\mu, \theta) \quad (8.3)$$

*We say that φ **dominates** θ **uniformly**.*

Proof.(*Proposition 5*). Let θ is Markov policy and

$$\theta(a|x) = \mathbb{P}^*\{a_t = a | x_{t-1} = x\} = \frac{\mathbb{P}^*\{x_{t-1} a_t = xa\}}{\mathbb{P}^*\{x_{t-1} = x\}} \quad (8.4)$$

$$(a \in A_t, \quad x \in X_{t-1}, \quad m+1 \leq t \leq n),$$

where \mathbb{P}^* is a probability measure in the space of ways L which is assigned to the initial distribution μ and to the policy π .

Remark 11. *The expression on the right side of (8.4) makes no sense for $\mathbb{P}^*\{x_{t-1} = x\} = 0$. So, for such x (in particular for killed states) we choose the arbitrary distribution on $A(x)$ instead of $\theta(\cdot|x)$.*

Let \mathbb{Q}^* denotes a probability distribution on space L which is assigned to the initial distribution μ and to the killed Markov policy θ .

The distribution \mathbb{Q}^* does not match with \mathbb{P}^* in the general case, but it is enough for proving (8.2) if any of $x_m, a_{m+1}, \dots, a_n, x_n$ and $x_{m+1}^*, x_{m+2}^*, \dots, x_n^*$ has the same probability distribution according to measures \mathbb{P}^* and \mathbb{Q}^* .

The following assertion holds

$$\begin{aligned} \omega(\mu, \pi) &= \sum_{t=m+1}^n \mathbb{P}^* q(a_t) + \sum_{t=m+1}^n \mathbb{P}^* c(x_t^*) + \mathbb{P}^* r(x_n), \\ \omega(\mu, \theta) &= \sum_{t=m+1}^n \mathbb{Q}^* q(a_t) + \sum_{t=m+1}^n \mathbb{Q}^* c(x_t^*) + \mathbb{Q}^* r(x_n). \end{aligned}$$

We shall use the mathematical induction to prove this.

The **basis** of induction: (8.2) holds for x_m because $\mathbb{P}^* = \mathbb{Q}^* = \mu$.

The **induction hypothesis**: let (8.2) holds for x_{t-1} . Let's check it for a_t .

Since θ is a killed Markov policy then

$$\mathbb{Q}^*\{x_{t-1}a_t = xa\} = \mathbb{Q}^*\{x_{t-1} = x\}\theta(a|x), \quad (a \in A_t, \quad x \in X_{t-1}). \quad (8.5)$$

Hence, from (8.4) and (8.5) we get

$$\begin{aligned} \mathbb{P}^*\{a_t = a\} &= \sum_{x \in X_{t-1}} \mathbb{P}^*\{x_{t-1}a_t = xa\} = \sum_{x \in X_{t-1}} \mathbb{P}^*\{x_{t-1} = x\}\theta(a|x) = \\ &= \sum_{x \in X_{t-1}} \mathbb{Q}^*\{x_{t-1} = x\}\theta(a|x) = \sum_{x \in X_{t-1}} \mathbb{Q}^*\{x_{t-1}a_t = xa\} = \mathbb{Q}^*\{a_t = a\}. \end{aligned}$$

So, our proposition holds for a_t .

The **induction hypothesis**: let (8.2) holds for a_t . Let show it for x_t .

From the definition of the transition function we get

$$\mathbb{P}^*\{a_t x_t = ax\} = \mathbb{P}^*\{a_t = a\}p(x|a), \quad (8.6)$$

$$\mathbb{Q}^*\{a_t x_t = ax\} = \mathbb{Q}^*\{a_t = a\}p(x|a). \quad (8.7)$$

From (8.6) and (8.7) it follows

$$\begin{aligned} \mathbb{P}^*\{x_t = x\} &= \sum_{a \in A_t} \mathbb{P}^*\{a_t x_t = ax\} = \sum_{a \in A_t} \mathbb{P}^*\{a_t = a\}p(x|a) = \\ &= \sum_{a \in A_t} \mathbb{Q}^*\{a_t = a\}p(x|a) = \sum_{a \in A_t} \mathbb{Q}^*\{a_t x_t = ax\} = \mathbb{Q}^*\{x_t = x\}, \quad (x \in X_t). \end{aligned}$$

Proposition 5 is proved.

Proof. (Proposition 6.) For proving this proposition we need the following lemma.

Lemma 1. Let f is a arbitrary function and let ν is a arbitrary probability distribution on countable space E .

If $\nu f < +\infty$ then the set $\Gamma = \{x : f(x) \geq \nu f\}$ has a positive measure ν , namely

$$\nu(\Gamma) > 0$$

(See proof in [1]).

According to (4.2) the condition (8.3) is equal to

$$\omega(x, \varphi) \geq \omega(x, \theta), \quad \forall x \in X_m.$$

Let separate the killed Markov policy θ by a product of the policies $\theta = \gamma\theta'$ where γ is the restriction of θ on X_m and θ' is the restriction of θ on $X_{m+1} \cup X_{m+2} \dots \cup X_n$.

According to the fundamental equation (5.1) it holds

$$\omega(x, \theta) = \gamma_x f,$$

where $\gamma_x(\cdot) = \gamma(\cdot|x)$ is the probability distribution on $A(x)$,

and $f(a) = q(a) + \omega'(p_a, \theta')$, $(a \in A_{m+1})$.

Since **Lemma 1** for $\tilde{A}(x) \subset A(x)$ it follows $\gamma_x(\tilde{A}(x)) > 0$, where $\tilde{A}(x) = \{a : f(a) \geq \gamma_x f = \omega(x, \theta)\}$. As a result, $\tilde{A}(x)$ is nonempty. If $\psi(x)$ is an arbitrary point of $\tilde{A}(x)$ then $f(\psi(x)) \geq \omega(x, \theta)$. But since the fundamental equation (5.1) we get $f(\psi(x)) = \omega(x, \psi\theta')$ and

$$\omega(x, \psi\theta') \geq \omega(x, \theta).$$

Let assume that condition (8.3) holds for the derived model \tilde{Z}^* . Then exists a simple policy φ' in \tilde{Z}^* which uniformly dominates the killed Markov policy θ' . According to the fundamental equation (5.1) and our assumption we get

$$\omega(x, \psi\varphi') = q(\psi(x)) + \omega'(p_{\psi(x)}, \varphi') \geq q(\psi(x)) + \omega'(p_{\psi(x)}, \theta') = \omega(x, \psi\theta') \geq \omega(x, \theta).$$

In the model Z^* simple policy $\varphi = \psi\varphi'$ dominates θ uniformly. Finally, (8.3) holds for model Z^* too.

Proposition 6. is proved.

9. Markovian property. Let $0 < k < n$, let use the killed policy ρ on the interval $[0, k]$ and killed policy π on the interval $[k, n]$. Doing analogically to **Definition 15** we can say that policy $\rho\pi$ is used.

Proposition 7. Let L_0 is the space of ways on the interval $[0, n]$, let L_k is the space of ways on the interval $[k, n]$ and let $P_x^{*\rho\pi}$ is the probability distribution which is assigned to the initial state x and to the killed policy $\rho\pi$, and analogically $P_y^{*\pi}$ is the probability distribution on L_k .

Then $\forall \xi = \xi(x_k a_{k+1} \dots x_n)$ on L_k holds

$$E_x^{*\rho\pi} \xi = E_x^{*\rho} [E_{x_k}^{*\pi} \xi]. \quad (9.1)$$

Proof. $\forall l = y_0 b_1 \dots b_k y_k b_{k+1} \dots y_n$ according to (3.1)

$$P_x^{*\rho\pi}(y_0 b_1 \dots y_n) = P_x^{*\rho}(c y_k) P_{y_k}^{*\pi}(y_k d), \quad (9.2)$$

where $c = y_0 b_1 \dots b_k$, $d = b_{k+1} \dots y_n$. Any function ξ on the space L_k can be interpreted on L_0 like function which does not depend on $x_0 a_1, \dots, a_k$. That is why we multiply the both sides of (9.2) by $\xi(y_k d)$ and sum up over all ways

$$E_x^{*\rho\pi}\xi = \sum_{cy_k} P_x^{*\rho}(cy_k) \sum_d P_{y_k}^{*\pi}(y_k d) \xi(y_k d). \quad (9.3)$$

But $P_{y_k}^{*\pi}(yd) = 0$ for $y \neq y_k$ and it follows

$$\sum_d P_{y_k}^{*\pi}(y_k d) \xi(y_k d) = \sum_{yd} P_{y_k}^{*\pi}(yd) \xi(yd) = F(y_k). \quad (9.4)$$

By substitution in (9.3) the expression from (9.4) and according to $\sum_{cy_k} P_x^{*\rho}(cy_k) F(y_k) = E_x^{*\rho} F(x_k)$, we get (9.1). **Proposition 7 is proved.**

Corollary 1.(Markovian property) Let $\nu(y) = P_\mu^{*\rho}\{x_k = y\}$ ($y \in X_k$) then $\forall \mu$

$$E_\mu^{*\rho\pi}\xi = E_\mu^{*\rho}[E_{x_k}^{*\pi}\xi].$$

In particular

$$E_\mu^{*\rho\pi}\xi(x_k a_{k+1} \dots x_n) = E_\nu^{*\pi}\xi(x_k a_{k+1} \dots x_n), \quad (9.5)$$

It follows from (9.1) and $\sum_{y \in X_k} \nu(y) P_y^{*\pi}\xi = E_\nu^{*\pi}\xi$.

The formula (9.5) shows that the probability distribution for a part of the trajectory does not depend on the distribution μ and policy ρ on the interval $[k, n]$. Namely, the probability forecast of the "future" (ξ) depends not on the "past" (μ, ρ), but only on the "present" (ν). Actually, it is already the **Markovain property**.

Let use the Markovian property for the assessment of a killed policy $\rho\pi$ on the intervals $[0, k]$ and $[k, n]$. Instead of ξ we take $\xi = \sum_{t=k+1}^n [q(a_t) + c(x_t^*)] + r(x_n)$ and by substituting in (9.5) we get

$$\omega(\mu, \rho\pi) = \sum_{t=1}^k E_\mu^{*\rho\pi}[q(a_t) + c(x_t^*)] + \omega(\nu, \pi) = \sum_{t=1}^k E_\mu^{*\rho}[q(a_t) + c(x_t^*)] + \omega(\nu, \pi). \quad (9.6)$$

The summation in (8.6) expresses the assessment $\omega(\mu, \rho)$ of policy ρ for a zero terminal reward, namely, $\omega(\mu, \rho\pi) = \omega(\mu, \rho) + \omega(\nu, \pi)$.

There is also another interpretation of (9.6). According to (4.2) and $\nu(y) = P_\mu^{*\rho}\{x_k = y\}$ ($y \in X_k$) we get

$$\omega(\nu, \pi) = \sum_y \nu(y) \omega(y, \pi) = E_\mu^{*\rho} \omega(x_k, \pi),$$

$$\omega(\mu, \rho\pi) = E_\mu^{*\rho} \left[\sum_{t=1}^k q(a_t) + \omega(x_k, \pi) \right]. \quad (9.7)$$

Hence, the assessment of killed policy $\rho\pi$ is equal to the assessment of the killed policy ρ with the terminal reward $\omega(\cdot, \pi)$ at the moment of time k .

10. Dynamic programming principle. Let Z^* be the model on the interval $[0, n]$ and let $0 \leq s < t \leq n$. Let $Z_{s,t}^*[f]$ denotes the model which is taken from the model Z^* by restriction of the interval $[0, n]$ to $[s, t]$. We define the terminal reward f at the moment of time t . Moreover, denote $\nu_s^t[f]$ as the assessment of the model Z_s^{*t} with the terminal reward f . Obviously, $\nu_s^t[f] = (VU)^{t-s} f = T^{t-s} f$ on X .

Since $\forall t \in [0, n]$ it holds

$$\nu_0^n[r] = \nu_0^t[\nu_t^n[r]] \text{ on } X_0 \text{ (} r \text{ on } X_n). \quad (10.1)$$

The equation (10.1) is equivalent to the optimality equations (7.1) and the condition $\nu^n = r$. It is called the ***Dynamic programming principle*** and it means that for the optimization of the decision on the interval $[0, n]$ with terminal reward r we must first optimize the decision on interval $[t, n]$ (with such terminal reward) and then optimize the decision on the interval $[0, t]$ with terminal reward $\nu_t^n[r]$.

In particular according to (9.1) it follows that if π'' is a uniform ε -optimal killed policy for Z_t^{*n} with terminal reward r and π' is a uniform ε -optimal policy for Z_0^{*t} with the terminal reward $\nu_t^n[r]$ then the killed policy $\pi = \pi''\pi'$ has the assessment $\nu_0^n[r]$ and is uniform ε -optimal for the model Z_0^{*n} (with terminal reward r).

References

- [1]. E.B. Dynkin, A.A. Yushkevich, *Markov Decision Processes*, M., (1975), 334 p. (Russian)
- [2]. E.A. Feinberg, A. Shwartz, *Introduction*, Handbook of Markov Decision Processes, Kluwer, (2002) (565 pages), pp.1-17. (English)
- [3]. A.G. Pakes, *Killing and Resurrection of Markov Processes*, Stochastic Models, V.13, I.2, (1997), pp.255-269. (English)
- [4]. R.E. Bellman, *Dynamic Programming*, Izdatelstvo inostrannoj literatury, (1960), 400 p. (Russian)

Nestor R. Parolya, Yaroslav I. Yeleyko

Ivan Franko National University of Lviv
 1, Universytetska str., 79000, Lviv, Ukraine
 Tel.: (8032) 239 45 31 (off.)